

Biocontainers 101

Using Containers in Bioinformatics

ITaP Research Computing Virtual Workshop Series

Yucheng Zhang, Lev Gorenstein

Feb. 25, 2022

Yucheng Zhang
Senior Life Science Scientist
ITaP Research Computing



Outline

- What are containers and why should we use them?
- Singularity basics
- Public container repositories
- Pull and use public biocontainers
- Deployed biocontainers on RCAC clusters
- Build your own biocontainers

What are containers?

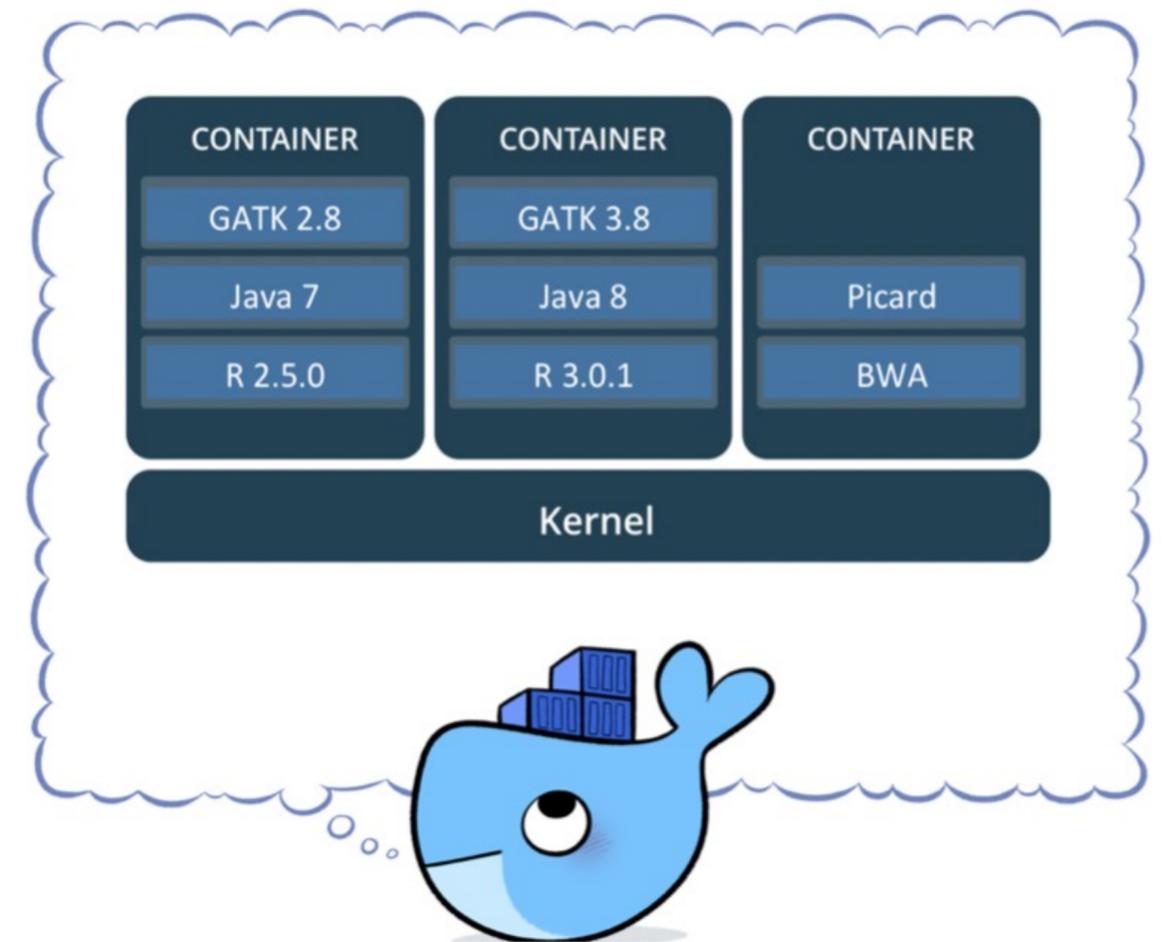
A **container** is an abstraction for a set of technologies that aim to solve the problem of how to get software to run reliably when moved from one computing environment to another.

A container **image** is simply a file (or collection of files) saved on disk that stores everything you need to run a target application or applications.

Registry is a place to store (and share) container images.

Why should we use containers?

- **Getting organized:** containers keep things organized by isolating programs and their dependencies inside containers.
- **Build once, run almost anywhere:** containers allow us to package up our complete software environment and ship it to numerous operating systems.
- **Reproducibility:** containers can ensure identical versions of apps, libraries, compilers, etc.



A real example

```
whatis("Pipeline for profiling microbial pathways from sequencing data ")
load("metaphlan2")
load("diamond")
load("MinPath")
load("RAPsearch")
load("usearch/8")
prepend_path("PATH", "/group/bioinfo/apps/apps/humann2-0.11.1/bin")
prepend_path("PYTHONPATH", "/group/bioinfo/apps/apps/humann2-0.11.1/lib/p
help([[
Notes:
HUMAN2: The HMP Unified Metabolic Analysis Network 2
version 0.11.1
Pipeline for profiling microbial pathways from sequencing data
http://huttenhower.sph.harvard.edu/humann2
Changes: /group/bioinfo/apps/apps/humann2-0.11.1/history.md

Includes full databases from:
http://huttenhower.sph.harvard.edu/humann2_data/chocophlan/full_chocophlan_p
http://huttenhower.sph.harvard.edu/humann2_data/uniprot/uniref_annotated/uni

Forum: https://groups.google.com/forum/#!forum/humann-users
Manual: http://huttenhower.sph.harvard.edu/humann2/manual
Tutorial: https://bitbucket.org/biobakery/biobakery/wiki/humann2

]])
```

REQUIREMENTS

1. **MetaPhlan 2.0**
2. **Bowtie2** (version ≥ 2.1) (see NOTE)
3. **Diamond** (0.9.0 > version $\geq 0.8.22$) (see NOTE)
4. **MinPath** (see NOTE)
5. **Python** (version ≥ 2.7)

```
diamond/0.6.12
diamond/0.7.1
diamond/0.7.8
diamond/0.7.9
diamond/0.8.36
diamond/0.9.6
diamond/0.9.26
```



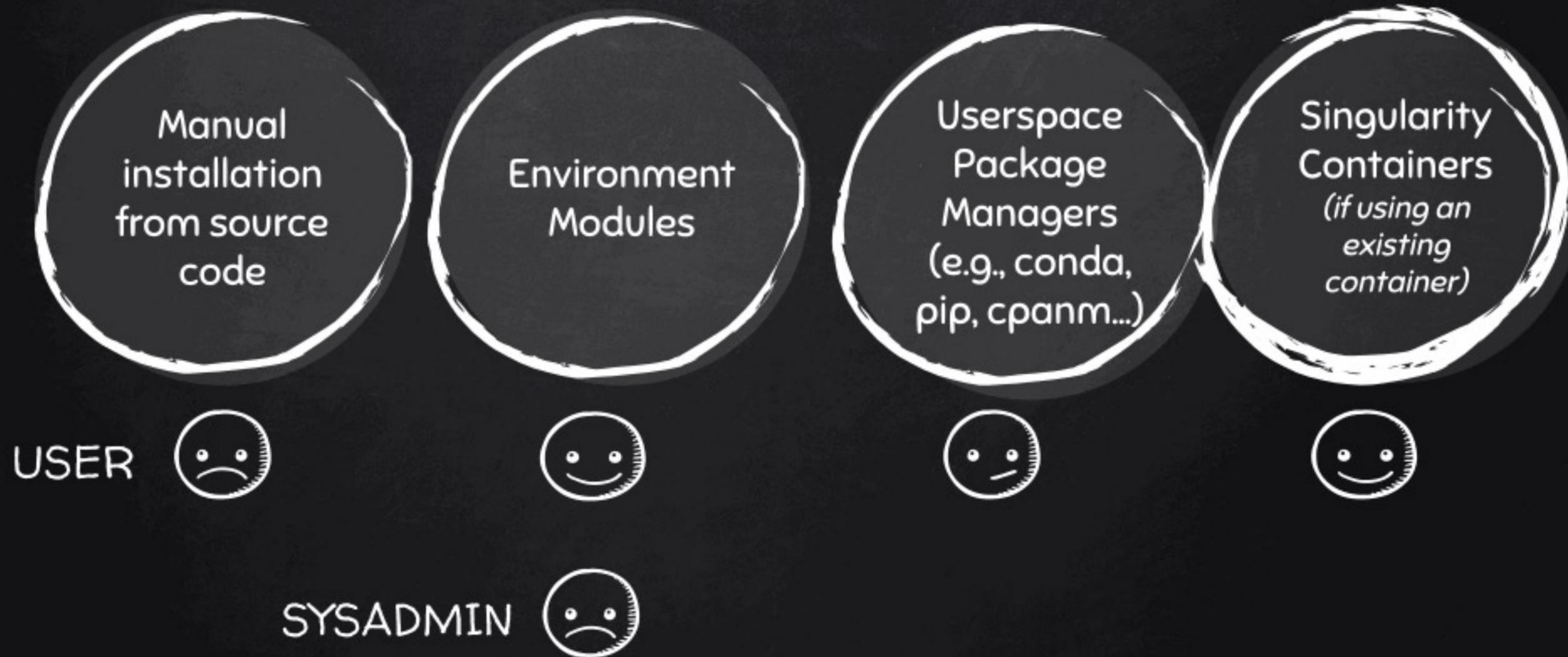
(D)

Why should we use containers on our clusters?

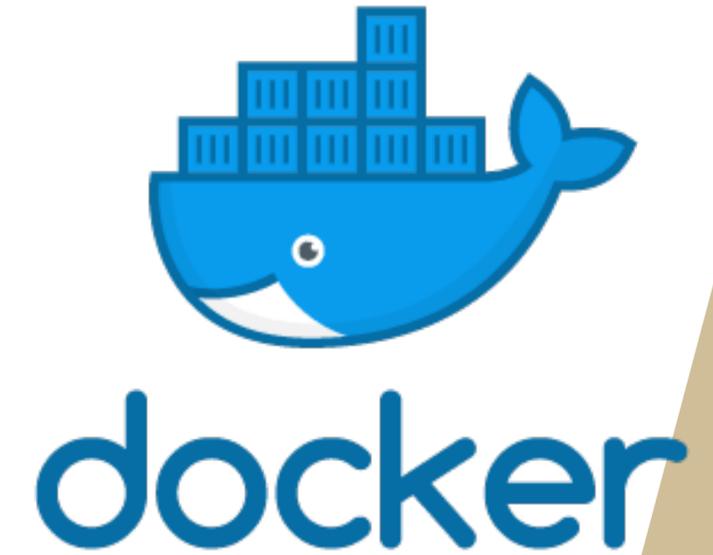
Enable you to install and use software easier

- ❖ Some software or packages has specific requirements for certain libraries such as **GLIBC**. Because our cluster's OS version is old, its libraries may not be compatible with your software. In such case, installing software into our clusters will be very challenging.
- ❖ However, if you build a container, you can get the latest everything and aren't limited by the cluster's OS version. **You are the master of your containers.**

INSTALLING SOFTWARE ON AN HPC CLUSTER, OVERSIMPLIFIED



Docker



The concept of containers emerged in 1970s, but they were not well known until the emergence of Docker containers in 2013.

Docker is an open source platform for building, deploying, and managing containerized applications.

Some concerns about the security of Docker containers on HPC: Docker gives superuser privileges, but we do not want users to have full, unrestricted admin/ root access.



Singularity

- ❖ Singularity was developed in 2015 as an open-source project by researchers at Lawrence Berkeley National Laboratory led by Gregory Kurtzer.
- ❖ Singularity is emerging as the containerization framework of choice in HPC environments.
 1. Enable researchers to package entire scientific workflows, libraries, and even data.
 2. Users do not need to ask their system admin (e.g., RCAC) to install software for them.
 - 3. Can use docker images.**
 4. Secure!
 - 5. Does not require root privileges.**



Singularity basics

Detailed singularity user guide is available at: sylabs.io/guides/3.8/user-guide

The main singularity command

```
singularity [options] <subcommand> [subcommand options ...]
```

- ❖ Build
- ❖ Pull
- ❖ Shell
- ❖ Exec





Singularity workflow on HPC

- 1 (Optional). **Build** singularity containers on a computer system where you have root or sudo privilege, e.g., your personal computer with singularity installed.
2. **Pull** the public containers or **transfer** your own containers to HPC.
3. **Run** singularity containers on the HPC system.

singularity pull

Download a container from a given URI.

singularity pull [output file] <URI>

Supported URIs include:

- ❖ **Library:** pull an image from singularity library
library://<user>/<collection>/<image>[:tag]
- ❖ **Docker hub:** pull an image from Docker Hub.
docker://<repository>/<image>[:tag]
- ❖ **Quay.io:** pull an image from Quay.io registry
docker://quay.io/<repository>/<image>[:tag]
- ❖ **http, https:** pull an image using the http(s?) protocol
e.g., https://library.sylabs.io/v1/imagefile/library/default/alpine:latest



Three useful image registries

1. Docker Hub (<https://hub.docker.org>)

- ❖ Online repository of Docker container images.
- ❖ As of Feb. 17, 2022, 8,842,825 available container images.

2. BioContainers (<https://biocontainers.pro/registry>)

- ❖ A community-driven project for bioinformatics containers.
- ❖ 10.6K tools, 45.4K versions, 222.5K containers and packages.
- ❖ The [Bioconda package index](#) lists all software available.
- ❖ The [Biocontainers registry](#) provides a searchable interface.

3. GALAXY project (<https://depot.galaxyproject.org/singularity/>)

- ❖ The BioContainers community also stored each singularity image in Galaxy depot.
- ❖ Can be pulled or ran using the HTTP protocol.



singularity pull example

singularity pull [options] name_to_save.sif URI

Let's pull bowtie2 from three different resources.

- **Docker hub** (<https://hub.docker.com/r/biocontainers/bowtie2/tags>)

```
singularity pull bowtie.2.4.1.sif docker://biocontainers/bowtie2:v2.4.1_cv1
```

- **Bioconda package index** (<https://bioconda.github.io/recipes/bowtie2/README.html>)

```
singularity pull bowtie.2.4.5.sif docker://quay.io/biocontainers/bowtie2: 2.4.5--py37hafa4d4c_1
```

- **GALAXY project** (<https://depot.galaxyproject.org/singularity>)

```
singularity pull bowtie.2.4.5.sif https://depot.galaxyproject.org/singularity/bowtie2:2.4.5--py39hbb4e92a_0
```

Recommendation: add the *--disable-cache* option to prevent image layers from being cached in `${HOME}/.singularity/cache`

```
singularity pull --disable-cache name_to_save.sif URI
```



singularity shell

Go inside the container and start an interactive shell

singularity shell myimage.sif

```
zhan4429@brown-fe00:~ $ cat /etc/*release
CentOS Linux release 7.7.1908 (Core)
NAME="CentOS Linux"
VERSION="7 (Core)"
ID="centos"
ID_LIKE="rhel fedora"
VERSION_ID="7"
PRETTY_NAME="CentOS Linux 7 (Core)"
ANSI_COLOR="0;31"
CPE_NAME="cpe:/o:centos:centos:7"
HOME_URL="https://www.centos.org/"
BUG_REPORT_URL="https://bugs.centos.org/"
```

```
zhan4429@brown-fe00:~ $ singularity shell ubuntu_latest.sif
Singularity> cat /etc/*release
DISTRIB_ID=Ubuntu
DISTRIB_RELEASE=20.04
DISTRIB_CODENAME=focal
DISTRIB_DESCRIPTION="Ubuntu 20.04.3 LTS"
NAME="Ubuntu"
VERSION="20.04.3 LTS (Focal Fossa)"
ID=ubuntu
ID_LIKE=debian
PRETTY_NAME="Ubuntu 20.04.3 LTS"
VERSION_ID="20.04"
HOME_URL="https://www.ubuntu.com/"
SUPPORT_URL="https://help.ubuntu.com/"
BUG_REPORT_URL="https://bugs.launchpad.net/ubuntu/"
PRIVACY_POLICY_URL="https://www.ubuntu.com/legal/terms-and-policies/privacy-policy"
VERSION_CODENAME=focal
UBUNTU_CODENAME=focal
Singularity> █
```

Type “**exit**” in the interactive shell to go back to host system



Bind mounts

- ❖ Programs running inside a container will not have access to directories and files outside of your home and the current directory.
- ❖ Singularity allows you to map directories on your host system to directories within your container using bind mounts.

```
singularity shell --bind hostdir1:containerdir1 --bind hostdir2:containerdir2 myimage.sif
```

Singularity binds several directories into the container image automatically. **\$HOME**, **/tmp** and **\$PWD** is the default list.

We also configured singularity to bind **/apps**, **/depot**, and **/scratch** on our clusters.



singularity exec

Run a command within a container

singularity exec myimage.sif **command**

For example:

```
singularity exec blast.2.11.0.sif blastx -query input.fasta -db swissprot -out blast.out
```

--bind option is also very useful for singularity exec

For example:

```
singularity exec --bind $HOME/data/:/data/ blast.2.11.0.sif blastx -query /data/input.fasta -db nr
```

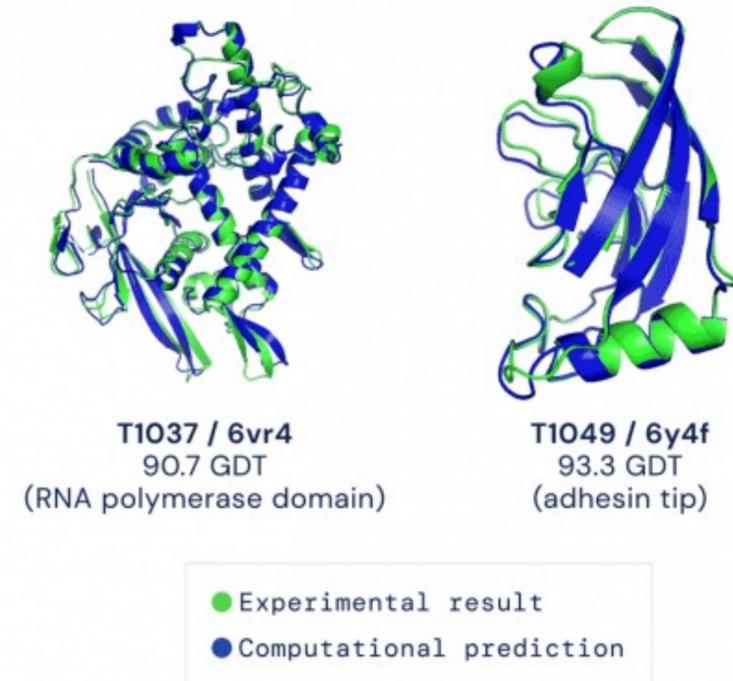
input.fasta is located in the host directory \$HOME/data/

AlphaFold

Deployed in all clusters, support both CPU and GPU.

\$ module load biocontainers

\$ module load alphafold/2.1.1



The full database (~2.2TB) has been downloaded and setup for users.

Usage:

```
run_alphafold.sh --flagfile=$AlphaDB --fasta_paths=XX --output_dir=XX ...
```

\$AlphaDB (/depot/itap/datasets/alphafold/full_db.ff) is a configuration file passed to AlphaFold containing the location of the database. Typically it should not be edited. Users can add other parameters based on your needs.

<https://github.com/deepmind/alphafold>



HOMER

Software for motif discovery and next-gen sequencing analysis

\$ module load biocontainers

\$ module load homer/4.11

Selected database have been downloaded for users.

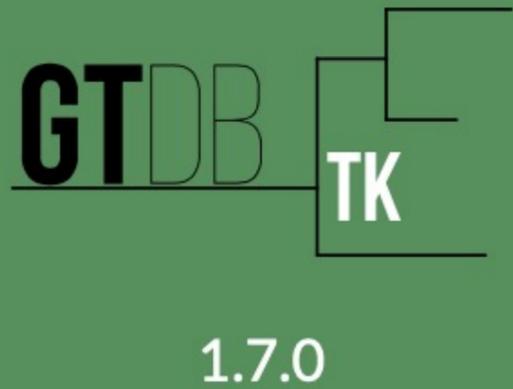
ORGANISMS: yeast, worm, mouse, arabidopsis, zebrafish, rat, human and fly.

PROMOTERS: yeast, worm, mouse, arabidopsis, zebrafish, rat, human and fly.

GENOMES: hg19, hg38, mm10, ce11, dm6, rn6, danRer11, tair10, and sacCer3.

Check installed databases:

\$ configureHomer.pl -list



Toolkit for assigning objective taxonomic classifications to bacterial and archaeal genomes based on the Genome Database Taxonomy [GTDB](#).

```
$ module load biocontainers
```

```
$ module load gtdbtk/1.7.0
```

GTDB-Tk reference data (R202) has been downloaded for users.

Example usage:

```
$ gtdbtk identify --genome_dir genomes --out_dir identify --extension gz --cpus 8
```

```
$ gtdbtk align --identify_dir identify --out_dir align --cpus 8
```

```
$ gtdbtk classify --genome_dir genomes --align_dir align --out_dir classify --extension gz --cpus 8
```

HUMAnN 3

Quantify species' contributions to community function

```
$ module load biocontainers
```

```
$ module load humann/3.0.0
```

Full **ChocoPhlAn**, **UniRef90**, **EC-filtered UniRef90**, **UniRef50**, **EC-filtered UniRef50**, and **utility_mapping** databases have been downloaded for users.

Check the database and config by:

```
$ humann_config --print
```

```
HUMAnN Configuration ( Section : Name = Value )
```

```
database_folders : nucleotide = /depot/itap/datasets/humann/chocophlan
```

```
database_folders : protein = /depot/itap/datasets/humann/uniref
```

```
database_folders : utility_mapping = /depot/itap/datasets/humann/utility_mapping
```

Run_dbcan

Automated CAZyme annotation

```
$ module load biocontainers
```

```
$ module load run_dbcan/3.0.2
```

Latest version of database has been downloaded and setup, including **CAZyDB.09242021.fa**, **dbCAN-HMMdb-V10.txt**, **tcdb.fa**, **tf-1.hmm**, **tf-2.hmm**, and **stp.hmm**.

Usage:

```
$ run_dbcan protein.faa protein --out_dir test1_dbcan
```

```
$ run_dbcan genome.fasta prok --out_dir test2_dbcan
```

https://github.com/linnabrown/run_dbcan



R-RNAseq

Customized R container for RNAseq analysis.

```
$ module load biocontainers
```

```
$ module load r-rnaseq/4.1.1-1
```

```
OR $ module load r-rnaseq/4.1.1-1-rstudio
```

Commands:

1. R
2. Rscript
3. rstudio (only exist in rstudio version)

BiocManager	1.30.16	readr	2.0.2
ComplexHeatmap	2.9.4	readxl	1.3.1
DESeq2	1.34.0	purrr	0.3.4
edgeR	3.36.0	dplyr	1.0.7
DEXSeq	1.40.0	stringr	1.4.0
pheatmap	1.0.12	forcats	0.5.1
limma	3.48.3	ggplot2	3.3.5
tibble	3.1.5	openxlsx	4.2.5
tidyr	1.1.4		

Thanks to Lev Gorenstein's r/4.1.1 base image, users can also install other packages, same with non-containerized R.



R-scRNAseq

Customized R container for scRNAseq analysis.

```
$module load biocontainers
```

```
$module load r-scrnaseq/4.1.1-1
```

```
OR $module load r-scrnaseq/4.1.1-1-rstudio
```

Commands:

1. R
2. Rscript
3. rstudio (only exist in rstudio version)

BiocManager 1.30.16	schex 1.8.0	tidyr 1.1.4
Seurat 4.1.0	CoGAPS 3.14.0	readr 2.0.2
SeuratObject 4.0.4	celldex 1.4.0	readxl 1.3.1
SeuratWrappers 0.3.0	dittoSeq 1.6.0	purrr 0.3.4
monocle3 1.0.0	DropletUtils 1.14.2	dplyr 1.0.7
SingleCellExperiment 1.16.0	miQC 1.2.0	stringr 1.4.0
scDbFinder 1.8.0	Nebulosa 1.4.0	forcats 0.5.1
SingleR 1.8.1	tricycle 1.2.0	ggplot2 3.3.5
scCATCH 3.0	pheatmap 1.0.12	openxlsx 4.2.5
scMappR 1.0.7	limma 3.48.3, 3.50.0	
rliger 1.0.0	tibble 3.1.5	

Build your own containers with singularity

The first step is to install singularity on your personal computer.

We have singularity version **3.8.0** on the cluster. To guarantee compatibility, please be sure to follow the installation guide for version 3.8 on your system (https://sylabs.io/guides/3.8/user-guide/quick_start.html).

```
$ sudo singularity build image.sif image.def
```

- ❖ Need to build using a computer with elevated privileges, then copy to cluster.
- ❖ If no access to such a computer, can also build in the cloud.



Remote builder

If you need to build an image from a system where you don't have admin privileges, we can build remotely using the [Sylabs Remote Builder](#).

To remotely build an image using singularity, go through the following steps:

1. Go to: <https://cloud.sylabs.io/>, and create a Sylabs account.
2. Create a new “Access Token”, and copy it to clipboard.
3. Login to our clusters, and run ``singularity remote login`` in terminal and paste the access token at the prompt.
4. Then you can remotely build your own singularity image on the cluster.

```
singularity build -r myimage.sif myimage.def
```

```
or singularity build --remote myimage.sif myimage.def
```

Once finished, the image will be downloaded automatically so that it's ready to use.

Singularity definition file

A **definition** file, or **def** file, is a recipe to build a container image with singularity. It is divided into two parts:

- 1. Header:** the Header describes the core operating system to build within the container.
- 2. Section:** each section is defined by a % character followed by the name of the particular section. Different sections add different content or execute commands at different times during the build process.

Detailed instruction on how to prepare a definition file is available at https://sylabs.io/guides/latest/user-guide/definition_files.html.

```

BootStrap: docker
From: debian:buster-slim
  
```

Header

```

#####
%post
#####
apt-get -y update
apt-get install -y curl wget nano bzip2 less

# miniconda
mkdir -p /opt
wget https://repo.continuum.io/miniconda/Miniconda3-latest-Linux-x86_64.sh
bash Miniconda3-latest-Linux-x86_64.sh -f -b -p /opt/conda
. /opt/conda/etc/profile.d/conda.sh
conda activate base
conda update --yes --all
conda config --add channels bioconda
conda config --add channels conda-forge
conda create -n prokka prokka==1.14.6

# create bind points for NIH HPC environment
mkdir /gpfs /spin1 /data /scratch /fdb /lscratch /vf
for i in $(seq 1 20); do ln -s /gpfs/gsfs$i /gs$i; done

# clean up
apt-get clean
conda clean --yes --all
  
```

Section

```

#####
%environment
#####
export LC_ALL=C
export PATH=/opt/conda/envs/prokka/bin:$PATH
  
```

def file for prokka 1.14.6 prepared by NIH HPC staff



Preferred bootstrap agents

1. **library:** images hosted in [Sylabs Cloud Library](#)
2. **docker:** images hosted in [Docker Hub](#)
3. **localimage:** images saved on your machine

Information about more bootstrap agents can be found in [Singularity user guide](#).

Bootstrap: docker

From: ubuntu:20.04

%labels

Author "Yucheng Zhang <zhan4429@purdue.edu>"

Version v0.935

%help

Singularity container with ANGSD v0.935. This container also installed htlib and samtools.

%post

update the system and install building essentials

apt-get -y update

apt-get -y install --no-install-recommends --no-install-suggests libssl-dev libcurl4-gnutls-dev libbz2-dev \
liblzma-dev libz-dev samtools gcc g++ git ca-certificates build-essential make zip wget unzip locales locales-all

clean up

apt-get -y autoremove && apt-get clean. && rm -rf /var/lib/apt/lists/*

Install htlib

SRC=/usr/local/src

mkdir -p \$SRC && cd \$SRC

git clone --recursive https://github.com/samtools/htlib.git

cd htlib && make

Install angsd

cd \$SRC && git clone https://github.com/ANGSD/angsd.git

cd angsd && make HTSSRC=\$SRC/htlib

#Symbolic link

chmod +x \$SRC/angsd/misc/realSFS

cd /usr/local/bin

ln -s \$SRC/angsd/angsd . && ln -s \$SRC/angsd/misc/realSFS .

ANGSD

Program for analyzing NGS data

Bootstrap: docker

From: continuumio/miniconda3

%labels

Author "Yucheng Zhang <zhan4429@purdue.edu>"

Version 2.4.3

%help

This container contains the latest version (v2.4.3) of aTRAM.

%post

```
conda install git
```

```
cd /opt/ && git clone https://github.com/juliema/aTRAM.git
```

```
cd aTRAM && chmod +x *.py
```

```
conda install python=3 numpy biopython psutil
```

```
conda install -c bioconda blast velvet trinity abyss spades exonerate
```

%environment

```
export PATH=/opt/aTRAM/:$PATH
```

aTRAM

automated target restricted
assembly method



Bootstrap: localimage

From: r-base:4.1.1.sif

%post

```
Rscript -e "install.packages('tidyverse')"
```

```
Rscript -e "install.packages('openxlsx', dependencies = TRUE)"
```

```
Rscript -e 'if (!requireNamespace("BiocManager", quietly = TRUE))
```

```
  install.packages("BiocManager")' \
```

```
&& Rscript -e
```

```
'BiocManager::install(c("limma", "edgeR", "DESeq2", "ComplexHeatmap", "DEXSeq"))'
```

If users want to build your own R containers, welcome to use our r-base images and recipes that are stored in </depot/itap/biocontainers/recipes/>



Bootstrap: localimage

From: r_4.1.1_rstudio.sif

%post

update the system and install building essentials

apt-get -y update

apt-get install -y gdal-bin libgdal-dev libudunits2-dev **## here you install required libraries**

clean up

apt-get -y autoremove && apt-get clean && rm -rf /var/lib/apt/lists/*

Seurat3

Rscript -e "install.packages('Seurat')"

monocle3

Rscript -e 'if (!requireNamespace("BiocManager", quietly = TRUE))

install.packages("BiocManager")' \

&& Rscript -e "BiocManager::install(c('BiocGenerics', 'DelayedArray', 'DelayedMatrixStats', 'limma', 'S4Vectors', 'SingleCellExperiment', 'SummarizedExperiment', 'batchelor', 'Matrix.utils'))"

Rscript -e "devtools::install_github('cole-trapnell-lab/leidenbase')"

Rscript -e "devtools::install_github('cole-trapnell-lab/monocle3')"

Cluster user guide

[Singularity section](#) contains instructions for using Singularity on RCAC clusters.

[Biocontainer collection section](#) contains instructions and examples for running bioinformatic containers.

Email

rcac-help@purdue.edu is our email support address. Send us an email any time.

Coffee hour consultations

In response to COVID-19, we are temporarily switching all our Coffee Hour Consultations to online only (<https://www.rcac.purdue.edu/coffee>). We offer several slots (2:00 to 3:30pm) each afternoon (Monday to Thursday) for private one-on-one consultations or questions of up to 30 minutes.

Additional bioinformatic tools

Contact me (zhan4429@purdue.edu) or Lev (lev@purdue.edu), if you want additional software added into RCAC biocontainers.

Thank you!
Questions?